

Supplementary information

Non-neutral evolution of H3.3-encoding genes occurs without alterations in protein sequence

Brejnev M. Muhire¹, Matthew A. Booker^{1,2} and Michael Y. Tolstorukov^{1,2,*}

¹Department of Molecular Biology, Massachusetts General Hospital and Harvard Medical School, Boston, MA 02114

²current address: Dana-Farber Cancer Institute, Boston, MA 02215

*corresponding author: tolstorukov@molbio.mgh.harvard.edu

Supplementary figures

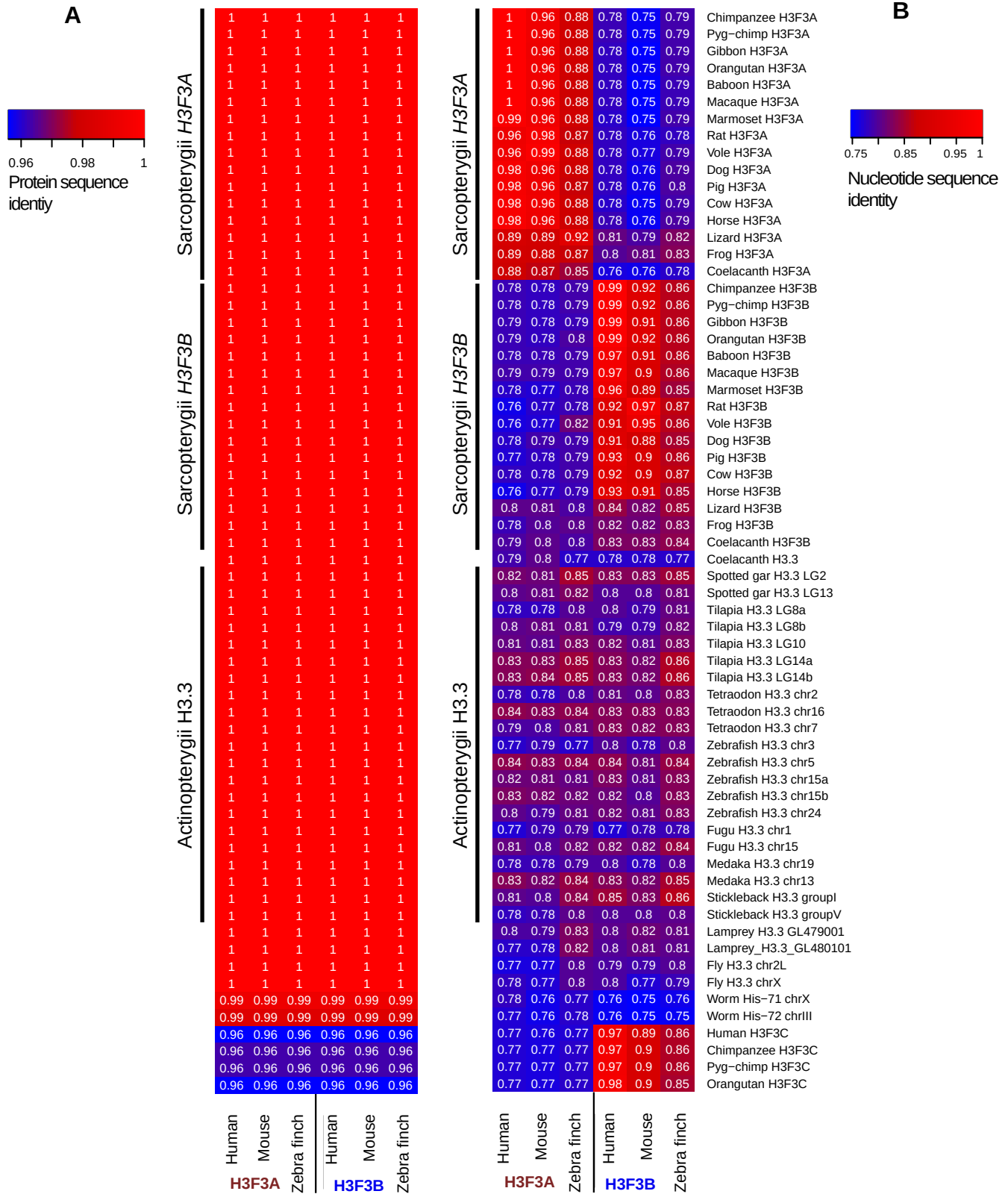


Figure S1. Sequence similarity of H3F3A, H3F3B and H3F3C in metazoa

A. Protein sequence identity computed between human, mouse and zebra finch H3.3 genes (H3F3A and H3F3B) and the H3.3 genes in analyzed metazoa genomes. The hominid-specific H3F3C is included at the bottom.

B. Similar identities as in (A) computed using coding DNA sequences.

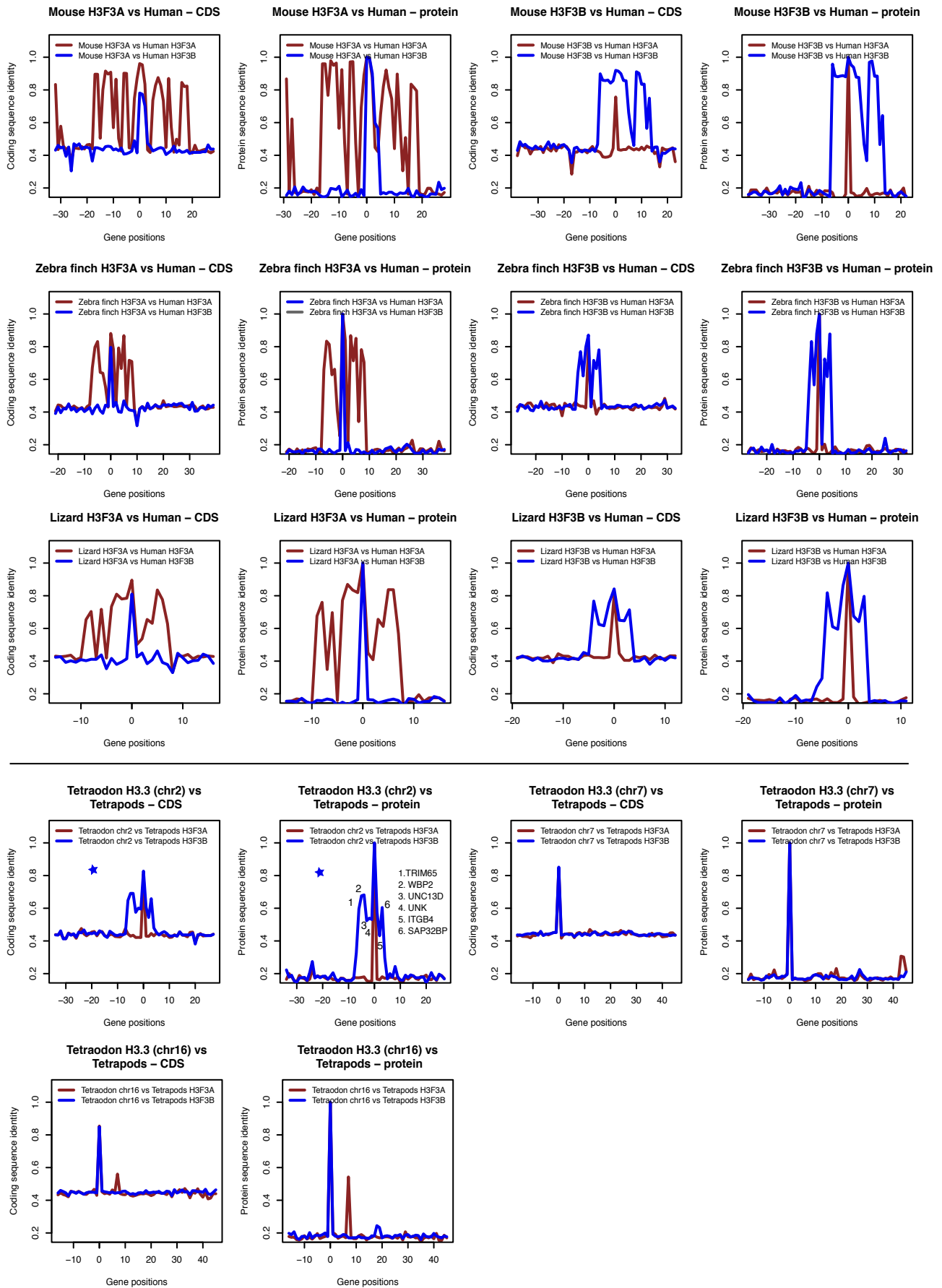


Figure S2A

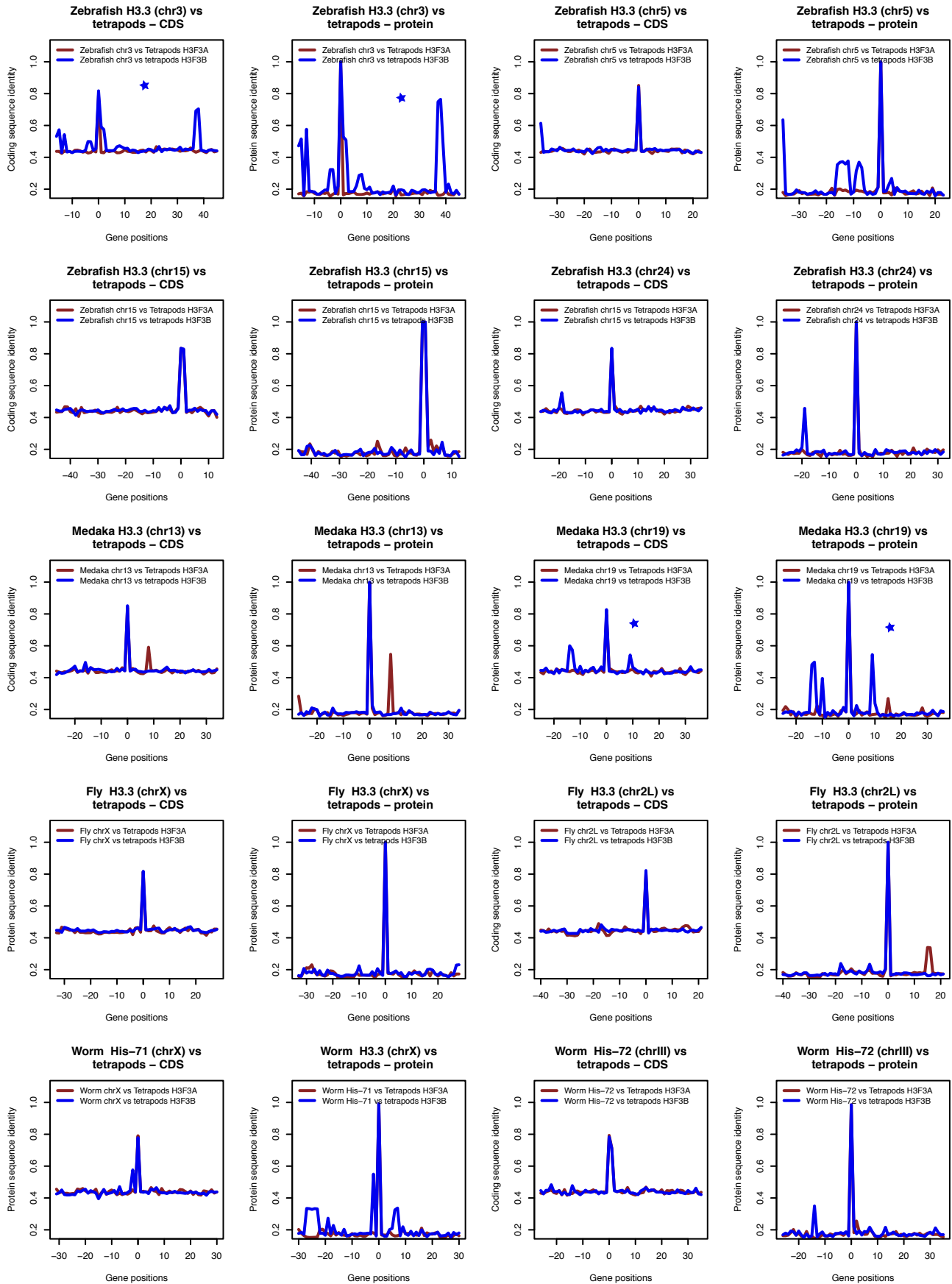


Figure S2B

Figure S2. Synteny around H3F3A and H3F3B genes

A, B Conservation of genes around H3.3 genes between human and other tetrapods (mouse, zebra finch and lizard; top three rows in A) and conservation genes around H3.3 genes between tetrapods and distant organisms (actinopterygian, fly and worm; the last two rows in A and all rows in B). Sequence conservation (identities) was measured for both nucleotide and protein sequences. The brown graph shows identities between genes around tetrapods H3.3 and human *H3F3A*, or genes conserved around non-tetrapod H3.3 and tetrapod *H3F3A* genes. Similarly, the blue graph shows identities between genes conserved around H3.3-encoding gene in a given organism and genes around Human *H3F3B* or tetrapod *H3F3B*. A blue star indicates non-tetrapod gene sharing syntenic genes with tetrapod *H3F3B*. The highest peaks at gene position 0 (100% protein identity for brown and blue lines) represent the identity between a given H3.3 gene and tetrapod *H3F3A* and *H3F3B*.

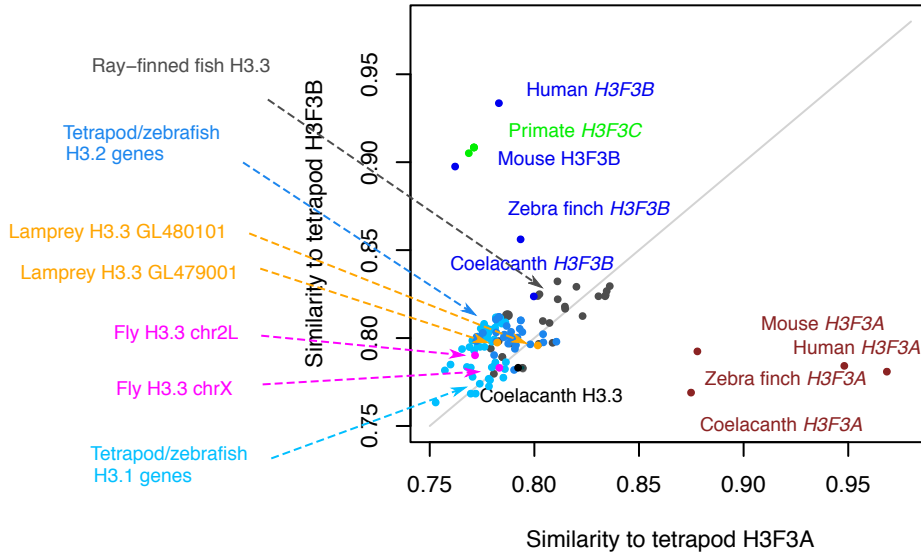


Figure S3. Comparison of tetrapod H3.3 genes to related genes in sarcopterygian and non-sarcopterygian lineages.

Average sequences similarity was estimated for the CDS of tetrapod (human, mouse, zebra finch) H3.3 genes (H3F3A, x-axis and H3F3B, y-axis) and CDS of each H3.3 gene in sarcopterygians, actinopterygians and more distant organisms (lamprey, fly). Additionally, CDS of tetrapod and zebrafish H3.1 and H3.2 genes were included in this analysis. Each point represents a gene and the organism name is written in the matching color. The sequence similarity represents percentage of the identical nucleotides in the sequence.

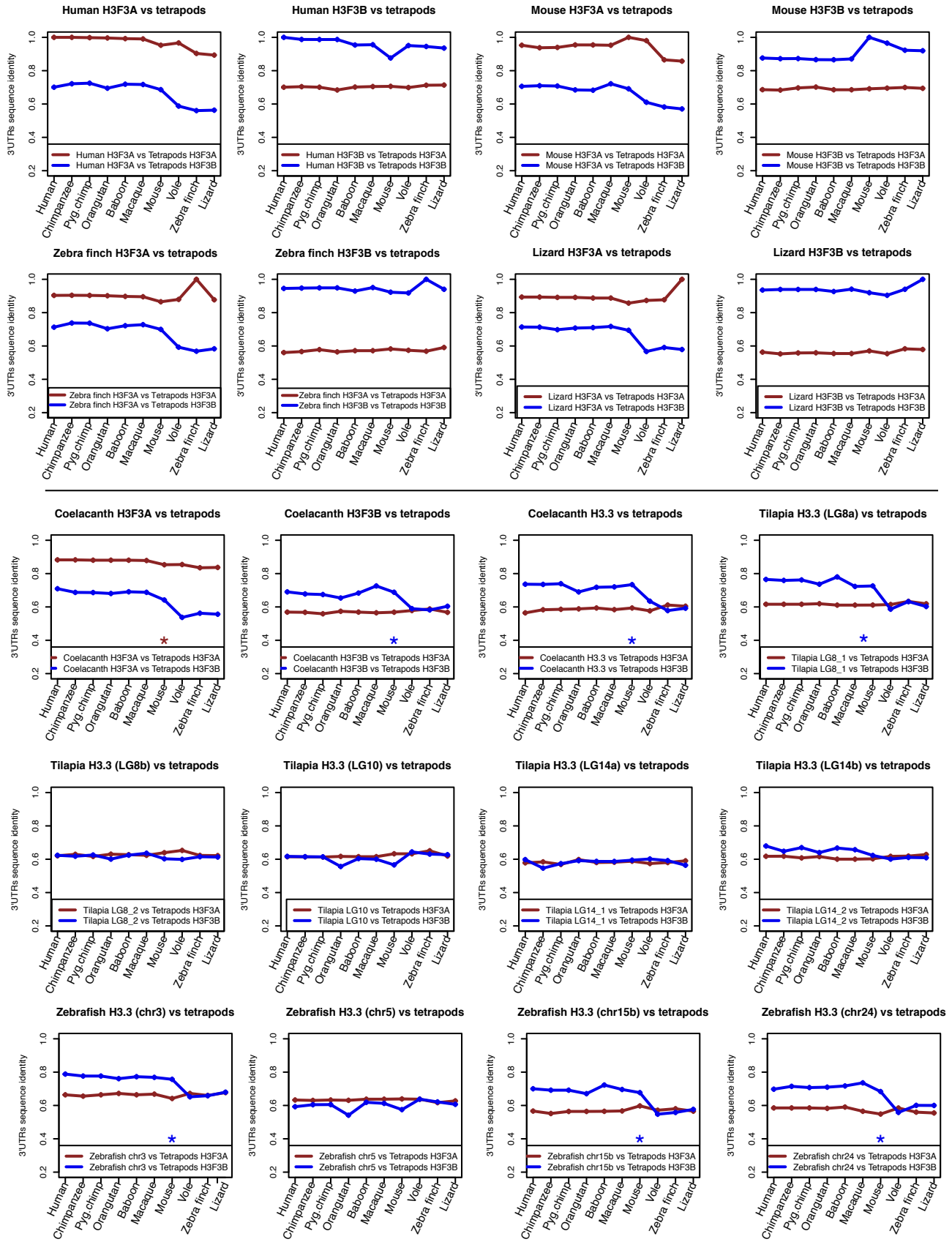


Figure S4A

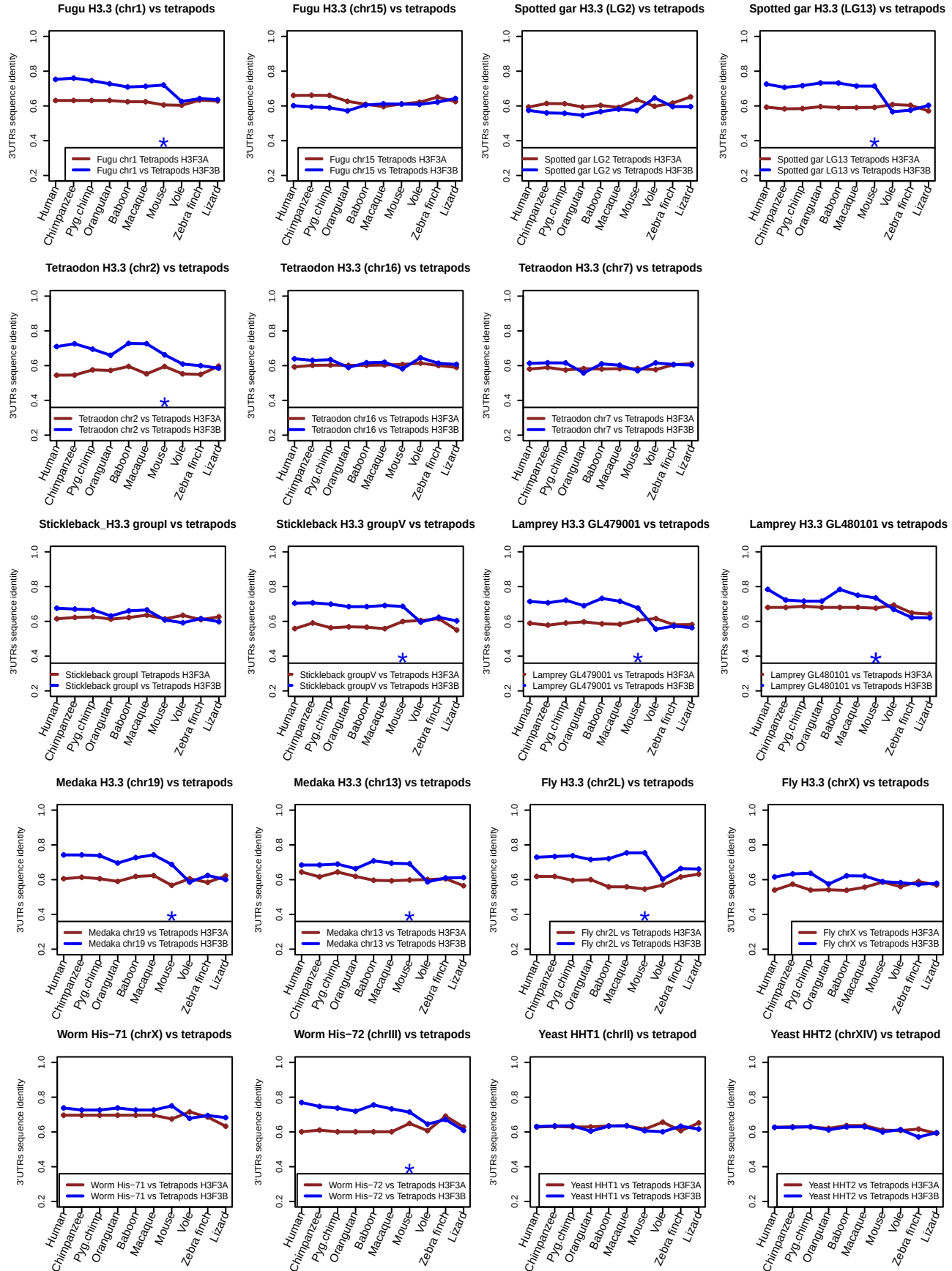


Figure S4B

Figure S4. Conservation of H3.3 genes 3'UTRs in metazoa

A, B H3.3 gene 3'UTRs sequence conservation among tetrapod organisms (top two rows in A) and H3.3 gene 3'UTRs conservation between tetrapods and non-tetrapod organisms (coelacanth, actinopterygians, fly and worm; last three rows in A and all rows in B). The blue line represents the sequence identity between a given H3.3 gene's UTR and the *H3F3B* 3'UTR from ten tetrapod organisms (x-axis). Similarly, the brown line shows the sequence identity between the 3' UTR of a given gene and the 3'UTR of tetrapod *H3F3A* genes. A blue asterisk indicates non-tetrapod H3.3 genes for which UTRs sequence is more similar to that of tetrapod *H3F3B* (blue) or tetrapod *H3F3A* (brown) in the majority tetrapod organisms included.

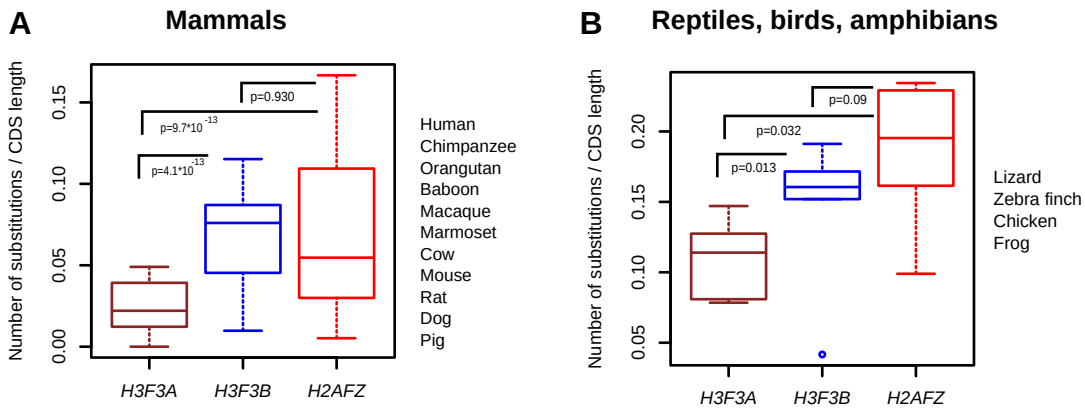


Figure S5. CDS conservation of histone variant genes in tetrapods

A, B. Pairwise nucleotide substitution scores (genetic distances) computed for two H3.3 genes (H3F3A, brown and H3F3B, blue), and H2AFZ gene (red) which was included in this analysis for comparison. The analysis was performed separately for mammalian (A) and other tetrapod genomes (reptiles, birds and amphibians; B). Comparison of scores was done with a Wilcox sum rank test and the organisms included are shown on the right side of each figure.

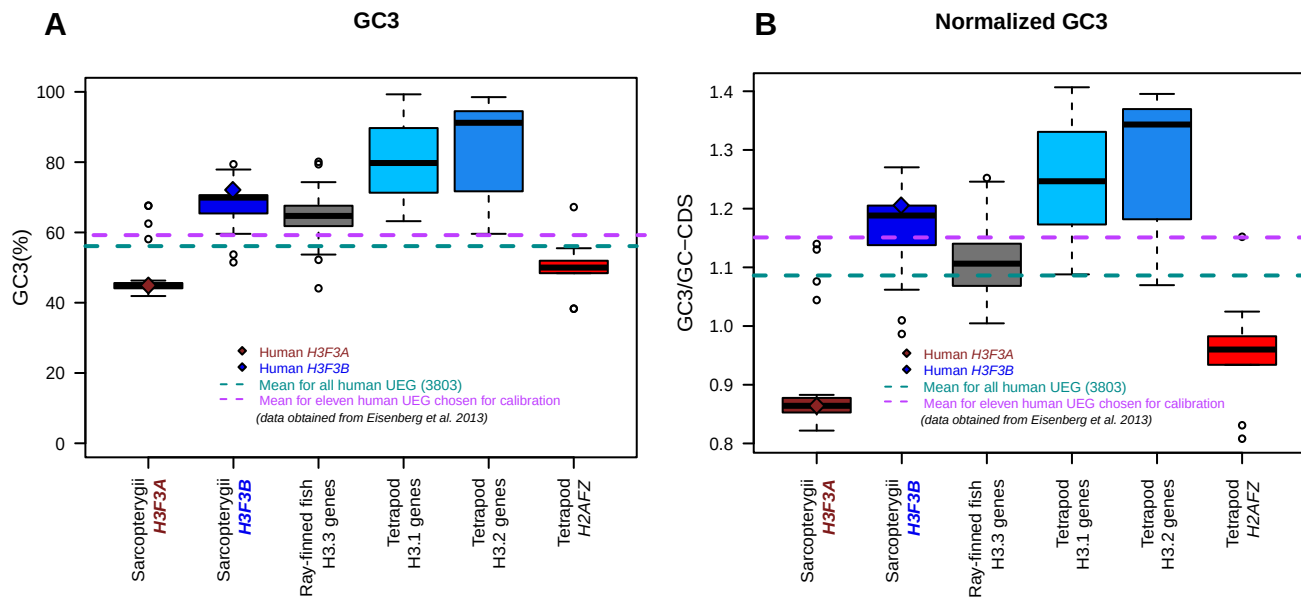
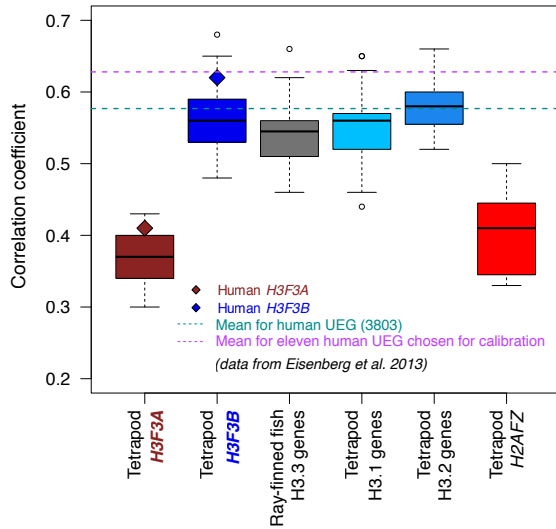


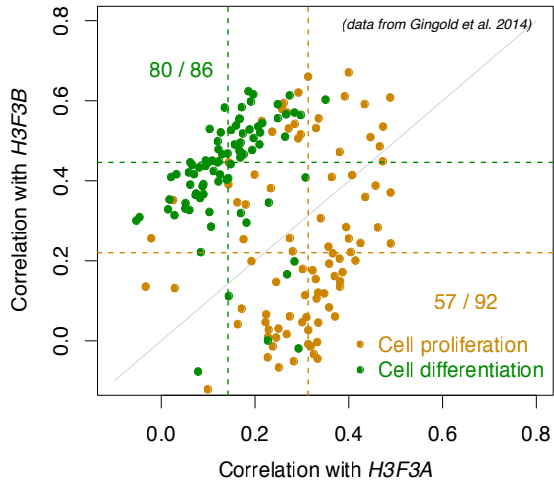
Figure S6. GC3 of H3.3 CDS.

A. Distribution of GC3 scores for tetrapod H3F3A and H3F3B, actinopterygian H3.3, tetrapod/zebrafish H3.1, H3.2 and H2AFZ genes. **B.** Distribution of normalized GC3 scores. For a particular gene, the normalized GC3 scores is the GC-content at codon 3rd position (GC3) divided by the GC-content of the whole gene (GC-CDS). The dashed lines represents average correlation computed for human ubiquitously expressed genes (UEG)⁴⁶. Cyan: the full set of UEGs; Magenta: a subset of eleven UEGs proposed because of their highly uniform and strong expression across human cell types.

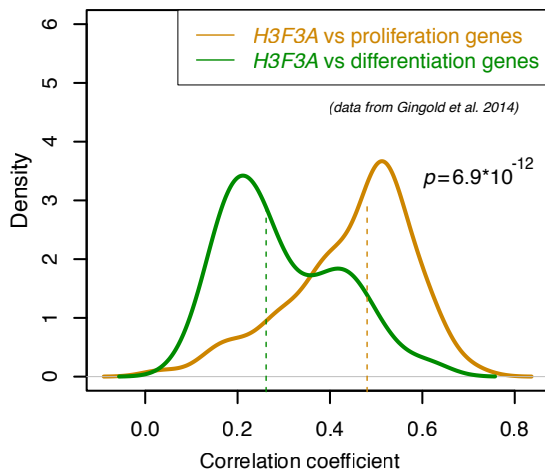
A Comparison with genomewide codon usage (codon frequencies)



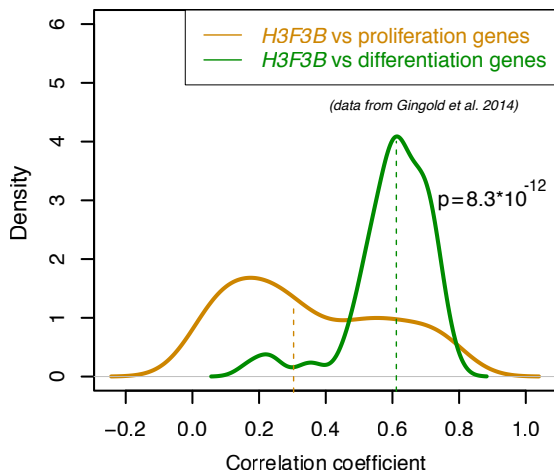
B Correlation of codon usages in H3.3 and proliferation- /differentiation-induced genes (codon frequencies)



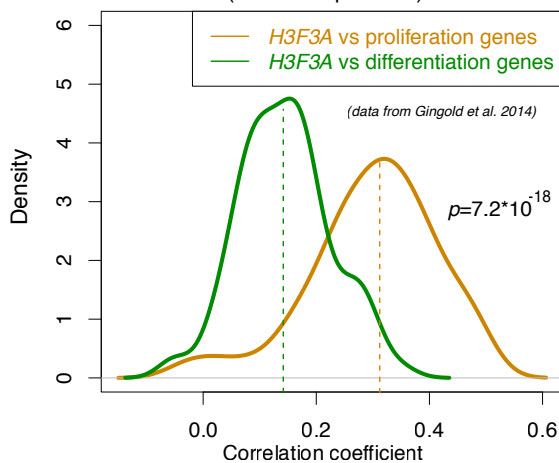
C Human *H3F3A* codon usage vs cell proliferation and cell differentiation (AA-specific codon frequencies)



D Human *H3F3B* codon usage vs cell proliferation and cell differentiation (AA-specific codon frequencies)



E Human *H3F3A* codon usage vs cell proliferation and cell differentiation (codon frequencies)



F Human *H3F3B* codon usage vs cell proliferation and cell differentiation (codon frequencies)

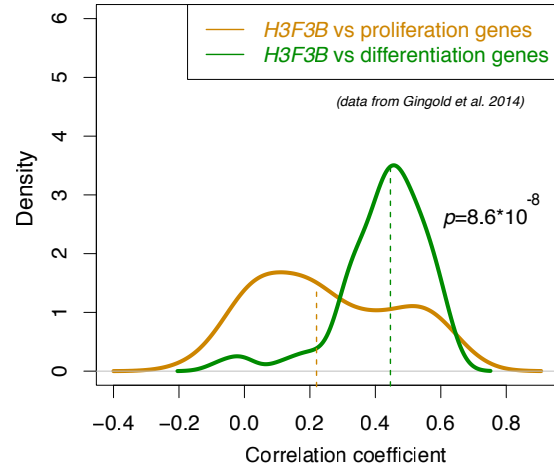


Figure S7

Figure S7. Distinct codon usage preferences in the H3.3 genes

A. Correlation between codon usage in the genes specified at x-axis and the genome-wide codon usage. The box plots represent the lineage distributions of the correlation coefficients calculated for the 'codon frequencies' of a corresponding gene with those estimated genome-wide (e.g. all tetrapod H3F3A genes vs. genome-wide frequencies). The brown and blue diamonds provide reference for human *H3F3A* and *H3F3B* respectively. The dashed lines represent average correlation computed for human ubiquitously expressed genes (UEG)⁴⁶. Cyan: the full set of UEGs; Magenta: a subset of eleven UEGs proposed because of their highly uniform and strong expression across human cell types. **B.** Correlation of human H3F3A and H3F3B 'codon frequencies' with those computed for the genes associated with cell proliferation (orange) and cell differentiation (green)⁴³. Each dot represents an individual gene from the corresponding group. The dotted lines indicate the correlation coefficient medians for each group and the H3.3 gene. **C,D.** Comparison of distributions of correlation scores of human H3F3A 'amino-acid specific codon frequencies' and those computed for genes associated with cell proliferation (orange), and cell differentiation (green)⁴³ (C). Similar comparisons performed for Human H3F3B and genes associated with cell proliferation (orange) and differentiation (green) (D). A vertical dashed line represents median of correlation scores obtained for a given H3.3 gene and a given group of genes. Comparison of scores was done with a Mann-Whitney test. **E,F.** Similar comparisons as in C and D respectively, performed based on 'codon frequencies'.

Supplementary table

Table S1. Number of H3.3 encoding genes per organism

No.	Name	Scientific name	Class / phylum	No. of H3.3
1	Human	<i>Homo sapiens</i>	Sarcopterygii	2
2	Chimpanzee	<i>Pan troglodytes</i>	Sarcopterygii	2
3	Pygmy chimpanzee	<i>Pan paniscus</i>	Sarcopterygii	2
4	Gorilla	<i>Gorilla gorilla</i>	Sarcopterygii	2
5	Orangutan	<i>Pongo abelii</i>	Sarcopterygii	2
6	Gibbon	<i>Nomascus leucogenys</i>	Sarcopterygii	2
7	Baboon	<i>Papio anubis</i>	Sarcopterygii	2
8	Macaque	<i>Macaca mulatta</i>	Sarcopterygii	2
9	Marmoset	<i>Callithrix jacchus</i>	Sarcopterygii	2
10	Mouse	<i>Mus musculus</i>	Sarcopterygii	2
11	Rat	<i>Rattus norvegicus</i>	Sarcopterygii	2
12	Prairie Vole	<i>Microtus ochrogaster</i>	Sarcopterygii	2
13	Opossum	<i>Monodelphis domestica</i>	Sarcopterygii	2
14	Dog	<i>Canis familiaris</i>	Sarcopterygii	2
15	Pig	<i>Sus scrofa</i>	Sarcopterygii	2
16	Cow	<i>Bos taurus</i>	Sarcopterygii	2
17	Horse	<i>Equus caballus</i>	Sarcopterygii	2
18	Chicken	<i>Gallus gallus</i>	Sarcopterygii	2
19	Zebra Finch	<i>Taeniopygia guttata</i>	Sarcopterygii	2
20	Lizard	<i>Anolis carolinensis</i>	Sarcopterygii	2
21	Frog	<i>Xenopus tropicalis</i>	Sarcopterygii	2
22	Coelacanth	<i>Latimeria chalumnae</i>	Sarcopterygii	3
23	Tilapia	<i>Oreochromis niloticus</i>	Actinopterygii	5
24	Zebrafish	<i>Danio rerio</i>	Actinopterygii	5
25	Fugu	<i>Takifugu rubripes</i>	Actinopterygii	3
26	Medaka	<i>Oryzias latipes</i>	Actinopterygii	3
27	Tetraodon	<i>Tetraodon nigroviridis</i>	Actinopterygii	3
28	Spotted gar	<i>Lepisosteus oculatus</i>	Actinopterygii	3
29	Stickleback	<i>Gasterosteus aculeatus</i>	Actinopterygii	3
30	Lamprey	<i>Petromyzon marinus</i>	Hyperoartia	2
31	Fruit fly	<i>Drosophila melanogaster</i>	Insecta	2
32	Worm	<i>Caenorhabditis elegans</i>	Nematoda	2